

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2000-99272
(P2000-99272A)

(43) 公開日 平成12年4月7日(2000.4.7)

(51) Int.Cl. ⁷	識別記号	F I	テーマコード(参考)
G 0 6 F 3/06	3 0 1 5 4 0	G 0 6 F 3/06	3 0 1 G 5 B 0 6 5 5 4 0

審査請求 未請求 請求項の数10 O L (全 14 頁)

(21) 出願番号 特願平10-272883

(22) 出願日 平成10年9月28日(1998.9.28)

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 日野 直樹

神奈川県小田原市国府津2880番地 株式会

社日立製作所ストレージシステム事業部内

(72) 発明者 中野 俊夫

神奈川県小田原市国府津2880番地 株式会

社日立製作所ストレージシステム事業部内

(74) 代理人 100068504

弁理士 小川 勝男

最終頁に続く

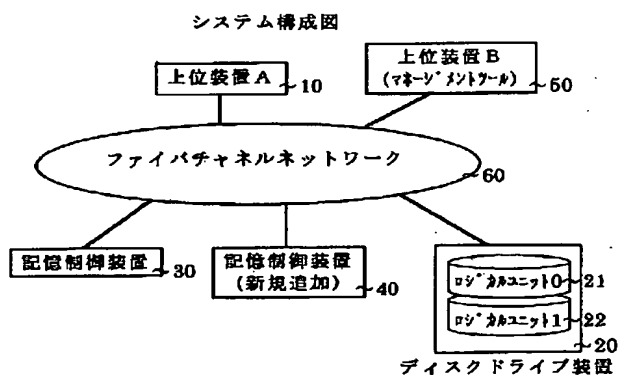
(54) 【発明の名称】 記憶制御装置及びこれを用いたデータ格納システムの取り扱い方法

(57) 【要約】

【課題】ファイバチャネルネットワーク上に新たな記憶制御装置をオンライン中に追加し、既に存在していた記憶制御装置からロジカルユニットの制御情報を引継ぐことで、これ以降の上位装置から当該ロジカルユニットへ行われる処理要求を新たな記憶制御装置が担当し、記憶制御装置間の負荷分散及び処理性能の向上を実現するファイバチャネル接続記憶制御装置を提供する。

【解決手段】ディスクドライブ装置20内の磁気ディスクドライブの構成情報やロジカルユニットの構成情報に代表されるロジカルユニット引継ぎ時に必要となる制御情報を記憶できる制御メモリを記憶制御装置30、40に持たせる。ファイバチャネルネットワーク上に新たな記憶制御装置40の追加時に、記憶制御装置30内の制御メモリの内容を記憶制御装置40の制御メモリにコピーする。これにより、記憶制御装置30、40間でロジカルユニットの引継ぎが可能であり、記憶制御装置間の負荷分散及び処理性能の向上を実現できる。

図 1



【特許請求の範囲】

【請求項 1】上位装置と、1つ又は複数の記憶装置との間に論理的に介在し、前記記憶装置を有する記憶制御装置において、
前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースに接続し、
前記記憶装置のロジカルユニットの制御情報を、前記記憶制御装置のメモリに格納し、前記上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置。

【請求項 2】上位装置と、1つ又は複数の記憶装置との間に論理的に介在し、前記記憶装置を有する記憶制御装置において、
前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースと、前記記憶装置と該記憶制御装置を論理的に結合するファイバチャネルにおける、前記記憶装置との接続に用いられるインタフェースとに接続し、
前記記憶装置のロジカルユニットの制御情報を、前記記憶制御装置のメモリに格納し、前記上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置。

【請求項 3】請求項 1 又は請求項 2 記載の記憶制御装置において、前記経路は、ANSI X3 T11 で標準化されたファイバチャネルであることを特徴とする記憶制御装置。

【請求項 4】1つ又は複数の記憶装置と、上位装置との間に論理的に介在し、該上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置であって、前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースに接続されるものと、前記上位装置及び前記記憶装置を有するデータ格納システムの取り扱い方法であって、
前記データ格納システムのオンライン中に、新たな記憶制御装置を追加するステップと、
前記新たな記憶制御装置の電源を投入するステップと、
前記上位装置で動作するツール又は前記記憶制御装置により、前記新たな記憶制御装置を特定するステップを有するデータ格納システムの取り扱い方法。

【請求項 5】1つ又は複数の記憶装置と、上位装置との間に論理的に介在し、該上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置であって、前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースと、前記記憶装置と該記憶制御装置を論理的に結合する経路における、前記記憶装置との接続に用いられるインタフェースとに接続されるものと、前記上位装置及び前記記憶装置を有するデータ格納

システムの取り扱い方法であって、

前記データ格納システムのオンライン中に、新たな記憶制御装置を追加するステップと、
前記経路の初期化処理を実行するステップと、
前記上位装置で動作するツール又は前記記憶制御装置により、前記新たな記憶制御装置を特定するステップを有するデータ格納システムの取り扱い方法。

【請求項 6】上位装置と、1つ又は複数の記憶装置との間に論理的に介在し、

10 前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースに接続し、

前記記憶装置のロジカルユニットを制御する為に用いる情報を、前記記憶制御装置のメモリに再設定されるまで保持し、前記上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置。

【請求項 7】請求項 4 又は請求項 5 記載のデータ格納システムの取り扱い方法において、

20 前記記憶制御装置は、更に、前記記憶装置のロジカルユニットの制御情報を格納する機能を有し、
更に、前記ツール又は前記記憶制御装置により、前記ロジカルユニットの制御情報の格納を開始するステップを有するデータ格納システムの取り扱い方法。

【請求項 8】1つ又は複数の記憶装置と、上位装置との間に論理的に既に介在し、該上位装置からの命令に応じて前記記憶装置に格納されたデータの授受を制御する記憶制御装置であって、前記上位装置と該記憶制御装置を論理的に結合する経路における、前記上位装置との接続に用いられるインタフェースに接続されるものと、前記上位装置及び前記記憶装置を有するデータ格納システムの取り扱い方法であって、

30 前記データ格納システムのオンライン中に、新たな記憶制御装置を追加するステップと、
前記経路の初期化処理を実行するステップと、
前記上位装置が前記既に介在する記憶制御装置を一意に識別するための情報を、前記新たな記憶制御装置が取得するステップを有するデータ格納システムの取り扱い方法。

【請求項 9】請求項 1、請求項 2 又は請求項 6 記載の記憶制御装置において、前記記憶装置は、磁気ディスク装置、光ディスク装置、光磁気ディスク装置、磁気テープ装置、又はライブラリ装置である記憶制御装置。

【請求項 10】請求項 4、請求項 5 又は請求項 8 記載のデータ格納システムの取り扱い方法において、前記記憶装置は、磁気ディスク装置、光ディスク装置、光磁気ディスク装置、磁気テープ装置、又はライブラリ装置であるデータ格納システムの取り扱い方法。

【発明の詳細な説明】

【0001】

50 【発明の属する技術分野】本発明は、磁気ディスク装

3

置、磁気テープ装置、光ディスク装置、光磁気ディスク装置、ライブラリ装置その他の記憶装置を制御し、これらと上位装置に対するデータの入出力（アクセス）を制御する記憶制御装置に関し、特に、記憶装置に対するアクセス経路としてファイバチャネルを用いた制御装置及びこれを用いたデータ格納システムに関する。

【0002】更に具体的には、本発明は、ANSI X3T11で標準化されたファイバチャネルを上位装置、磁気ディスク装置その他の記憶装置及び記憶制御装置のインターフェースとするコンピュータシステムに関し、コンピュータシステムのオンライン中に新たな記憶制御装置を追加し、記憶制御装置間の負荷の分散又は複数の記憶制御装置が行っていた機能の統合を可能とするデータ格納システムに関する。

【0003】

【従来の技術】記憶制御装置の増設に関しては特開平7-20994号公報に大型計算機の記憶システムにおいて、当該記憶システム内での上位装置との接続アダプタ、ディスクドライブ装置との接続アダプタ及び共有キャッシュメモリの増設による規模拡張及び活線挿抜機能について記載されており、前記記憶システム内の内部バス上に前記上位装置との接続アダプタ、ディスクドライブ装置との接続アダプタ及び共有キャッシュメモリを順次増設していく手法を提案している。

【0004】上記特開平7-20994号公報の技術は前記大型計算機の記憶システムで一般的に用いられている共有キャッシュメモリ方式を採用しており、前記上位装置との接続アダプタが内部バス上の共有キャッシュメモリに置かれたロジカルユニットの制御情報を逐次読込むことで、前記上位装置との接続アダプタ間でのロジカルユニットの制御情報の引継ぎ過程を必要とすることなく、任意のロジカルユニットへのアクセスを可能としている。

【0005】また、キャッシュメモリを個々に持つ記憶制御装置の多重化に関しては、特開平7-160432号公報に記憶制御装置の二重化を行い、片方の記憶制御装置を障害時のスタンバイ記憶制御装置として待機させておくことで記憶制御装置の冗長化を実現する手法について記載されている。さらに、特開平8-335144号公報に記憶制御装置の二重化による負荷分散及び記憶制御装置の障害発生時の冗長化について記載されている。

【0006】

【発明が解決しようとする課題】上記特開平7-20994号公報の技術では、ロジカルユニットの制御情報を内部バス上の共有キャッシュメモリ上に格納してしまうため、1記憶システム筐体内での規模拡張しか行なえない。

【0007】上記特開平7-160432号公報及び特開平8-335144号公報の技術では、二重化する記

4

憶制御装置において、予めオフライン時に各記憶制御装置を物理的に接続し、正常時及び障害時における各記憶制御装置が担当するロジカルユニットの設定を行なう必要がある、さらにロジカルユニットの制御情報を内部バスを用いて転送するため、初期設定に依存し、オンライン中の拡張性を持たない。

【0008】

【課題を解決するための手段】柔軟な機能拡張を行なうには、内部バスや共有キャッシュメモリの制限を受けないように、新たな記憶制御装置を追加するバス又は経路は上位装置との接続に使用されるインターフェースを用い、さらに追加する記憶制御装置が個々のキャッシュメモリを持ち、記憶装置の論理的まとまりであるロジカルユニット（1つの論理ドライブ又は複数の論理ドライブ群）の制御情報を前記キャッシュメモリへコピーすれば良い。

【0009】本発明では、ANSI X3T11で標準化されたファイバチャネルを上位装置、磁気ディスク装置その他の記憶装置及び記憶制御装置のインターフェースとするデータ格納システムを構築する。

【0010】そして該システムにおいて、ファイバチャネルネットワーク上に新たな記憶制御装置をオンライン中に追加し、既存の記憶制御装置から制御情報を引継ぎ、新たに追加した記憶制御装置と既存の記憶制御装置との間の負荷分散を行う。

【0011】新たな記憶制御装置は、コンピュータシステムのオンライン中に既存のファイバチャネルネットワークへ追加されると、ファイバチャネルネットワーク上の既存の記憶制御装置から、ロジカルユニット単位で制御情報を、ファイバチャネルネットワークを経由して取得する（引継ぐ）手段又は機能を有する。制御情報の取得又は引継ぎ完了により、所定のロジカルユニットが追加された記憶制御装置の配下に入る。

【0012】追加された新たな記憶制御装置は、その配下に入ったロジカルユニットに対して、上位装置からコマンド処理要求があれば、これを処理する。こうして新たに追加された記憶制御装置と、既存の記憶制御装置との間の負荷分散を行う。

【0013】本発明の負荷分散は、ファイバチャネルネットワーク上の既存の複数の記憶制御装置が、オンライン中の負荷分散及び統合を行う手段又は機能を持たない場合であっても、適用可能である。尚、本発明を実施する場合には、既存の記憶制御装置に、制御情報の引継ぎを行う機能を設定する必要がある。

【0014】更に、ロジカルユニット単位の制御情報引継ぎに際し、オペレータがファイバチャネルネットワーク、LAN（Local Area Network）その他のネットワークに接続された上位装置上で動作するマネージメントツールを用いて、又は、記憶制御装置のパネルを用いて、ロジカルユニットの引継ぎ開始や、引

継ぎ中の動作モードを指定する手段又は機能を上位装置又は記憶制御装置に設け、新たな記憶制御装置の追加を適切なタイミングで行うことができる。

【0015】更に具体的には、新たにファイバチャネルネットワークへ接続された記憶制御装置は、ロジカルユニット単位の制御情報の引継ぎに際し、上位装置から記憶制御装置を一意に識別する情報であるN_Portアドレスを、前記ファイバチャネルネットワーク上に接続され且つ引継ぎ対象のロジカルユニットを所有している既存の記憶制御装置から引継ぐ手段又は機能を有する。これにより、上位装置から発せられたコマンド処理要求の経路を変更しなくて済む。

【0016】本発明の記憶制御装置は、制御メモリを有し、当該制御メモリ上に、磁気ディスク装置その他の記憶装置の種別、記憶装置の記憶容量、ブロック数、各記憶装置の状態及びRAID (Redundant Arrays of Inexpensive Disks) 構成情報を格納する物理ドライブ制御テーブルやロジカルユニットの先頭LBA (Logical Block Address) 及び最後尾LBAを格納するロジカルユニット制御テーブルに代表されるロジカルユニット引継ぎの際に必要な制御情報を記憶する手段又は機能を有する。

【0017】そして、同一ファイバチャネルネットワーク上に新たに追加される記憶制御装置は前記ファイバチャネルネットワーク上に既に存在していた記憶制御装置上の制御メモリから前記制御情報を当該新たに追加される記憶制御装置の制御メモリ上にコピーする手段又は機能を有する。

【0018】記憶制御装置上の前記制御メモリは揮発性メモリであっても良いが、通常は、不揮発性メモリとする。また、制御メモリの内容(制御テーブル上の設定情報ほか)をファイバチャネルネットワークに接続された磁気ディスク装置その他の記憶装置へ書き込むことで、制御メモリを不揮発性としたのと同様の効果が得られる。つまり、万一、電源の瞬断が生じて、記憶制御装置は、当該設定情報の再設定が行われるまでは、恒久的に設定情報を維持することができる。

【0019】また、本発明のロジカルユニット単位で制御情報をオンライン中に引継ぐことを用いれば、複数の記憶制御装置が行なっていたいくつかの処理を、任意の記憶制御装置へ統合することもできる。

【0020】新たな記憶制御装置の追加の仕方については、1) 新たな記憶制御装置をファイバチャネルネットワークへ物理的に接続した後、2) 当該記憶制御装置側からファイバチャネルネットワークに対してリンクリセットを発行し、ファイバチャネルネットワーク内への論理的なログインを行う。3) その後、オペレータがファイバチャネルネットワークへ接続された上位装置上のマネージメントツールから、又は、記憶制御装置のパネル

から、新たに追加された記憶制御装置の認識を行い、

4) 新たに追加された記憶制御装置が担当するロジカルユニットの指定、ロジカルユニットの引継ぎ開始の指示及びファイバチャネルネットワークに既に存在する記憶制御装置に対して当該ロジカルユニットの引継ぎ中に上位装置から処理要求があった場合の応答方法を設定し、5) ファイバチャネルネットワーク上に既に存在していた記憶制御装置から新たに追加した記憶制御装置に対して、ロジカルユニットの制御情報の転送を行なわせる。

6) 更に、前記ファイバチャネルネットワーク上に既に存在していた記憶制御装置において、上位装置から記憶制御装置を一意に識別する情報であるN_Portアドレスを複数設定しておき、7) ファイバチャネルネットワークへ新たに追加された記憶制御装置が前記N_Portアドレスの一部を引継ぐことにより、上位装置からのコマンド処理要求経路の変更を不要とする。8) 更に、コンピュータシステムのオンライン中にファイバチャネルネットワーク上に存在する記憶制御装置が同一ファイバチャネルネットワークに接続された別の記憶制御装置からロジカルユニット単位で制御情報を引継ぎ、

9) これ以降の上位装置から当該引き継いだロジカルユニットに対して行われるコマンド処理要求を担当し、10) ファイバチャネルネットワーク上に存在する記憶制御装置間の負荷分散を行う。11) 必要に応じ、複数の記憶制御装置が行なっていたいくつかの処理を、任意の記憶制御装置へ統合することを、ロジカルユニットの制御情報をオンライン中に引継ぐ技術を適用して行う。

【0021】

【発明の実施の形態】 本発明の一実施例を図面を用いて以下に説明する。尚、記憶装置として磁気ディスク装置(以下、ディスクドライブ装置と記す)のみを用いた実施例としたが、光ディスク装置、テープ記憶装置であっても良い。

【0022】図1は本発明の一実施形態であり、ファイバチャネルネットワークにFC-AL (Fibre Channel Arbitrated Loop) の接続形態(トポロジ)を採用した場合のハードウェア構成図である。

【0023】図1において、10はデータ処理を行う中央処理装置(CPU)をもつ上位装置である。60はFC-AL・トポロジで動作し、FC-ALハブにて各種装置が接続されたファイバチャネルネットワークである。

【0024】20は上位装置A(10)からのデータを格納しておくディスクドライブ装置であり、複数の磁気ディスクドライブで構成される。ディスクドライブ装置20を構成する複数のディスクドライブを論理的に分割し、分割した区画を任意のRAIDレベルに定義することで、ディスクドライブ障害時における冗長性を持たせ、ディスクドライブ障害時におけるデータ消失を防ぐ

ことができる。この区画をRAIDグループと呼ぶ。このRAIDグループをさらに論理的に分割したSCSI (Small Computer System Interface) のアクセス単位である領域をロジカルユニットと呼び、この領域はLUN (Logical Unit Number) という番号をもつ。本実施例ではディスクドライブ装置20はロジカルユニット0 (21) 及びロジカルユニット1 (22) を持っているが、ロジカルユニットの個数は図1に示す2個でなくても良い (SCSI-3規格に準拠したファイバチャネルではファイバチャネルID毎に最大64個のロジカルユニットを持つことができる)。また、図1ではディスクドライブ装置が1個の場合の例を示しているが、ファイバチャネルネットワーク上にディスクドライブ装置が複数あってもかまわない。複数のディスクドライブ装置からなる大容量のデータ格納システムでは記憶制御装置の処理能力不足に陥ることが多々あり、本発明はコンピュータシステムが記憶制御装置の処理能力不足に陥った際に特に有効である。

【0025】30、40は上位装置A (10) とディスクドライブ装置20の間に論理的に介在し、上位装置A (10) とディスクドライブ装置20間の情報 (データ) の授受を制御する記憶制御装置である。特に、記憶制御装置30は本発明で必要となるファイバチャネル上に既に存在していた記憶制御装置を示し、記憶制御装置40は新たに追加される記憶制御装置を示す。図1では記憶制御装置30及び40が各々1装置となっているが、記憶制御装置30、40は複数あっても良い。また、図2において示される記憶制御装置30、40の内部構成図において、31はファイバチャネルネットワーク上でのデータ転送を制御し、さらに上位装置10から送られてくるコマンドの解析を行い、データをキャッシュ制御部34へDMA (Direct Memory Access) 転送するファイバチャネル制御部である。32は記憶制御装置の動作を制御するマイクロプログラム及びロジカルユニット引継ぎ時に必要となる制御情報を格納する不揮発性の制御メモリである。33は記憶制御装置全体を制御する中央処理装置 (CPU) である。34はキャッシュへのデータの読み書きを制御するキャッシュ制御部、35はディスクドライブ装置20への書き込みデータ及び読み出しデータを一時的に格納しておくキャッシュである。36は記憶制御装置の動作設定を変更及び参照するパネルである。

【0026】50は記憶制御装置30、40の動作を制御するマネージメントツールをもつ上位装置である。本実施例では上位装置B (50) にマネージメントツールを持たせているが、マネージメントツールは上位装置A (10) にあってもかまわない。さらに、マネージメントツールを搭載する上位装置B (50) と記憶制御装置30、40のインターフェースはLANその他の遠隔操

作が可能なネットワークシステムであっても良い。

【0027】次に、上位装置A (10) が記憶制御装置30経由でディスクドライブ装置20とデータ転送を行う場合を例にとり、制御の流れ、データの流れを説明する。上位装置A (10) がアクセス要求を出すと、その要求を認識したファイバチャネル制御部31はCPU33に対して割り込み要求を発行する。CPU33は、上位装置A (10) からのコマンドを解析し、制御メモリ32内にある磁気ディスクドライブの記憶容量、ブロック数、各磁気ディスクドライブの状態及びRAID構成情報を格納する物理ディスクドライブ制御テーブルやロジカルユニットの先頭LBA (Logical Block Address) 及び最後尾LBAを格納するロジカルユニット制御テーブルから情報を読み出す。

【0028】上位装置A (10) からのアクセス要求がライトコマンドの場合は、CPU33はファイバチャネル制御部31にデータ転送を指示し、上位装置A (10) から転送されたライトデータをキャッシュ制御部34を経由してキャッシュ35に格納し、上位装置A (10) に対して、ファイバチャネル制御部31がライト完了報告を行う。ライト完了報告後、CPU33はファイバチャネル制御部31を制御し、ファイバチャネルネットワークを介してディスクドライブ装置のロジカルユニット21または22へ前記ライトデータ及び冗長データを書き込む。この時、ライトデータを格納するロジカルユニットのRAIDレベルがRAID5であった場合、ライトデータを格納するために行われる旧データと旧冗長データの読み出し処理、新冗長データの生成処理及び前記ライトデータと新冗長データの格納処理という記憶制御装置30及びディスクドライブ装置20に対して非常に負荷の高い処理を必要とする。これをWriteペナルティと呼ぶ。このWriteペナルティ処理においては、ディスクドライブ装置20へのアクセスが多数発生することも性能劣化の一因となるが、それ以前にディスクドライブ装置20を制御する記憶制御装置30のマイクロプログラムの実行に多大な時間を要し、記憶制御装置30のCPU33が処理能力不足になることも多い。本発明ではオンライン中に記憶制御装置40を増設することで、記憶制御装置30の処理能力不足を解消する効果がある。

【0029】一方、上位装置A (10) からのアクセス要求がリードコマンドの場合は、CPU33は、ファイバチャネル制御部31に指示を出し、当該アクセス要求のデータブロックが格納されたディスクドライブ装置20内のロジカルユニット21、22へアクセスしてデータを読み出し、キャッシュ制御部34を経由してキャッシュ35へリードデータを格納する。キャッシュ35にリードデータを格納した後、CPU33はファイバチャネル制御部31に指示を出し、キャッシュ35に格納したリードデータを上位装置A (10) に転送し、転送終

了後、上位装置 A (10) ヘリード完了報告を行なう。

【0030】次にファイバチャネルネットワーク 60 の特長を説明する。ファイバチャネルは最大 10 km のデータ転送距離及び最大 100 MB/s のデータ転送速度を実現する長距離、高速インタフェースである。また、ファイバチャネルは上位論理層に位置する SCSI、IP、IPI 等の種々のプロトコルを下位論理層のファイバチャネルプロトコルにマッピングする機能を持っており、上位装置 A (10) からは SCSI や IP のように異なったプロトコルを持つデバイスを同一のファイバチャネルネットワークに接続し、データ転送を行なうことができる。即ち、他のインタフェースと論理的に互換性を持つ。

【0031】ファイバチャネルでは 3 つの接続形態 (トポロジ) が定義されている。一つ目は上位装置とデバイスが一对一に接続されるポイント・ツー・ポイント・トポロジである。二つ目は複数の上位装置とデバイスが一つのループを形成して接続される FC-AL・トポロジである。FC-AL 接続は FC-AL ハブと呼ばれる装置を介して接続される。三つ目はファブリック・トポロジであり、Fabric スイッチと呼ばれる装置を介して上位装置とデバイスがスター状に配置される。

【0032】図 1 で示す本実施例はファイバチャネルネットワークに FC-AL・トポロジを採用し、さらにファイバチャネルネットワーク上に接続された各種装置は SCSI マッピングプロトコルを用いて動作するコンピュータシステムであるが、図 3 で示すように Fabric スイッチを用いて各種装置を接続し、Fabric・トポロジで動作する構成も考えられる。また、図 1 においてはファイバチャネルネットワークが 1 個の FC-AL ループを形成しているが、図 4 で示すように複数の FC-AL ハブを用いてファイバチャネルネットワークを複数の FC-AL ループで形成し、さらに上位装置 A

(10) と記憶制御装置 30、40 間のファイバチャネルネットワークと記憶制御装置 30、40 とディスクドライブ装置 20 間のファイバチャネルネットワークを異なるループとし、ロジカルユニット制御情報の引継ぎにおいては上位装置 A (10) と記憶制御装置 30、40 間のファイバチャネルネットワークと記憶制御装置 30、40 とディスクドライブ装置 20 間のファイバチャネルネットワークの双方の引継ぎが必要なコンピュータシステムも考えられる。

【0033】ファイバチャネル上の情報のやりとりは Ordered Set と呼ばれる信号レベルの情報と、フレームと呼ばれる固定フォーマットを持った情報を用いて行われる。Ordered Set の代表的なものにはフレームの先頭を識別するために使用する SOF (Start Of Frame)、フレームの最後尾を識別するために使用する EOF (End Of Frame)、ループ上をフレームが転送されていない事を示す

Idle、FC-AL・ループ初期化開始要求に使用する LIP (Loop Initialization) 等の信号がある。

【0034】次にファイバチャネルがデータのやりとりを行なう基本単位であるフレームについて説明を行なう。フレームは機能に基づいてデータフレームとリンク制御フレームとに大別される。データフレームは情報を転送するために用い、データフィールドのペイロード部に SCSI 等の上位プロトコルで使用するデータ、コマンドを搭載する。

【0035】一方、リンク制御フレームは一般にフレーム配信の成功あるいは不成功を示すのに使われる。フレームを受領したことを示す ACK フレームやログインする場合に転送に関するパラメータを通知したりするフレーム等がある。

【0036】フレームのフォーマットについて、図 5 を用いて説明を行なう。フレーム 70 は、SOF 71、フレームヘッダ 72、データフィールド 73、CRC 74、EOF 75 で構成される。

【0037】SOF (Start Of Frame) 71 はフレームの先頭に置く 4 バイトの識別子である。EOF (End Of Frame) 75 はフレームの最後につける 4 バイトの識別子で、SOF 71 と EOF 75 によりフレームの境界を示す。

【0038】フレームヘッダ 72 はフレームタイプ、上位プロトコルタイプ、送信元と送信先の N_Port アドレス情報を含む。N_Port アドレスはファイバチャネルネットワークに接続される上位装置 10、50 及び記憶制御装置 30、40 等のアドレスを表わす。

【0039】データフィールド 73 の先頭部には上位レイヤのヘッダを置くことができる。これにデータそのものを運ぶペイロード部が続く。CRC (Cyclic Redundancy Check) 74 はフレームヘッダとデータフィールドのデータをチェックするための、4 バイトのチェックコードである。

【0040】次に図 6 にフレームヘッダのフォーマットを示す。フレームヘッダ 80 の D_ID (Destination ID) 81 はフレーム受け取り側の、また、S_ID (Source ID) 82 はフレーム送信側の N_Port アドレスである。FC-AL・トポロジでは 24 ビットの N_Port アドレスの内、下位 8 ビットを用いてアドレスを示し、特に FC-AL の N_Port アドレスを AL-PA (Arbitrated Loop Physical Address) と呼ぶ。

【0041】次にフレームを構成するデータフィールド 73 のペイロードの一部である FCP_CMND (Fibre Channel Protocol for SCSI Command) 及び FCP_RSP (Fibre Channel Protocol for SCSI

Response) について説明する。

【0042】FCP_CMNDのフォーマットを図7に示す。FCP_LUN (FCP Logical Unit Number) フィールド91にはコマンドを発行するロジカルユニット番号が指定される。FCP_CN TL (FCP Control) 92にはコマンド制御パラメータが指定される。

【0043】FCP_CDB (FCP Command Descriptor Block) 93にSCSIインターフェースでのコマンド授受に用いられるSCSI CDBが格納され、Operation Code 95ではInquiryやRead、Write、Mode sense、Mode Selectコマンド等を示す番号が格納される。さらにLUNや論理ブロックアドレス、転送ブロック長等の情報も格納される。FCP_DL (FCP Data Length) 94には本コマンドで転送されるデータ量がバイト数で指定される。このように構成されたフレームによってファイバチャネルネットワーク上でSCSIコマンドのやりとりが行われる。

【0044】FCP_RSPのフォーマットを図8に示す。FCP_RSP100はFCP_CMNDで指示されたSCSIコマンドの動作の結果を報告するために使用される。FCP_STATUS101にはコマンドが正常に完了したことを報告するGoodステータス、コマンドが異常終了したことを示すCheck Conditionステータス、デバイスが上位装置からのコマンドを受け付けられない状態にあることを示すBusyステータス等のSCSIステータスが格納される。

【0045】次に、FCP_CMND90におけるInquiry、Mode Sense、Mode Selectの役割について説明する。図9にInquiryコマンドのシーケンスを示す。Inquiryコマンドは上位装置10、50がファイバチャネルネットワークに接続された各種装置のデバイスタイプやサポートしている機能、装置のメーカー名、製品名を調べるために使用される。これらの情報はInquiryコマンドを受領した記憶制御装置30、40が上位装置10、50に対して送信するInquiryデータに格納されている。

【0046】図10にMode Senseコマンドのシーケンスを示す。Mode Senseコマンドは上位装置10、50がファイバチャネルネットワークに接続している各種装置のパラメータを参照するために使用される。Mode Senseコマンドで参照できるページは複数の固定ページとベンダユニークページが存在し、Page Codeにて参照するページを指定できる。これらの情報はMode Senseコマンドを受領した記憶制御装置が送信するMode Senseデータに格納されている。

【0047】一方、図11で示すMode Sele c

tコマンドは前記Mode Senseコマンドで参照できるパラメータを変更するためのコマンドである。パラメータを変更するには変更したいページをPage Codeにて指定し、Mode Selectパラメータリストを記憶制御装置に対して送信すれば良い。

【0048】次に、本発明において、記憶制御装置間でのロジカルユニットの引継ぎを行なう場合に必要となる情報を格納した制御メモリ32について説明を行なう。制御メモリ32は不揮発メモリであり、万一の電源瞬断時にも格納された情報を恒久的に維持することができる。

【0049】制御メモリ32には記憶制御装置30、40の動作を制御するマイクロプログラムの他に、物理ディスクドライブ制御テーブル110及びロジカルユニット制御テーブル120が格納されている。物理ディスクドライブ制御テーブルのフォーマットを図12で示す。物理ディスクドライブ制御テーブル110において、物理ディスクドライブ番号111はディスクドライブ装置内の磁気ディスクドライブに一意に割り当てられた番号である。物理ディスクドライブ位置112は当該磁気ディスクドライブの位置を示す論理的なアドレスを格納している。記憶容量113及びブロック数114は当該磁気ディスクドライブの全記憶容量及び全ブロック数を、RAIDグループ番号115は当該磁気ディスクドライブが所属するRAIDグループ番号を格納している。また、状態116は当該磁気ディスクドライブが使用可能なオンライン状態か、または使用不可能な閉塞状態かを示す情報が格納されている。また、記憶装置の種類117はファイバチャネルネットワーク60上の記憶装置が磁気ディスク装置、光ディスク装置、光磁気ディスク装置、磁気テープ装置、または各種ライブラリ装置のいずれであるかを識別するために用いる。

【0050】ロジカルユニット制御テーブルのフォーマットを図13で示す。ロジカルユニット制御テーブル120において、ロジカルユニット番号121はディスクドライブ装置内のロジカルユニットに対して一意に割り当てられる番号である。RAIDグループ番号は当該ロジカルユニットが所属するRAIDグループ番号であり、RAIDレベル123はそのRAIDレベルである。先頭アドレス124及び最終アドレス125は当該ロジカルユニットのRAIDグループでの位置を示すために使用する先頭ロジカルブロックアドレス及び最終のロジカルブロックアドレスである。前記制御メモリ32内の物理ディスクドライブ制御テーブル110及びロジカルユニット制御テーブル120の情報は他の装置が発行するMode Senseコマンドによる参照、及びMode Selectコマンドによるパラメータ変更が可能とする。

【0051】次に本発明における、ファイバチャネルネットワーク60への新たな記憶制御装置40のオンライ

10

20

30

40

50

ン中の追加及びロジカルユニットの引継ぎ処理手順を示す。この処理手順は記憶制御装置 40 の動作シーケンスと上位装置 B (50) に搭載されたマネージメントツールの動作シーケンスに大別される。図 14 に記憶制御装置 40 の動作シーケンス、図 15 に上位装置 B (50) に搭載されたマネージメントツールの動作シーケンスを示す。

【0052】まず、ファイバチャネルネットワーク 60 への記憶制御装置 40 の追加及び AL-PA 確定手順について説明する。ファイバチャネルネットワーク 60 を構成する FC-AL ハブと新たに追加する記憶制御装置 40 をファイバチャネルケーブルで接続した後、新たに追加する記憶制御装置 40 の電源を投入する (オン)。電源をオンにすることで、記憶制御装置 40 はリンク初期化処理を実行し、ファイバチャネルネットワーク 60 上をフレームの送受信可能な状態にする。

【0053】次に、新たに追加された記憶制御装置 40 はループ初期化処理を実行する。FC-AL・トポロジをもつファイバチャネルネットワークに新たに追加された記憶制御装置 40 は電源投入時には有効な AL-PA を持っていない。そこで、新たに追加された記憶制御装置 40 がループ初期化処理を実行することで有効な AL-PA が割り当てられる。AL-PA が確定した記憶制御装置 40 は一旦処理を停止し、上位装置 B (50) に搭載されたマネージメントツールからの指示を待つ。

【0054】次に、前記マネージメントツールの動作について説明を行なう。前記マネージメントツールは図 16 で示すマネージメントツール管理テーブル 130 をファイル形式で所有している。前記マネージメントツール管理テーブル 130 はファイバチャネルネットワーク 60 上に接続されている記憶制御装置 30 の AL-PA 131 及び担当 LUN 132 の情報が格納されている。

【0055】前記マネージメントツールはファイバチャネルネットワーク 60 へ新たに追加された記憶制御装置 40 の AL-PA を入手するために、FCP_CMND D90 の 1 種である Inquiry コマンドをファイバチャネルネットワーク 60 上に接続された全ての装置に対して発行する。そして、その応答である Inquiry データに格納された装置のメーカー名及び製品名を調べ、さらにフレーム送信側の AL-PA を参照し、マネージメントツール管理テーブル 130 に登録されていない装置を見つける。これにより、新規に追加された記憶制御装置 40 の AL-PA が特定でき、その AL-PA をマネージメントツール管理テーブル 130 に新規に書き込む。

【0056】新規に追加された記憶制御装置 40 の AL-PA が特定されると、マネージメントツールは記憶制御装置 40 に対して、引継ぎを行なうロジカルユニット、当該ロジカルユニットを担当していた記憶制御装置 30 の AL-PA 及び当該ロジカルユニットを担当して

いた記憶制御装置 30 のロジカルユニット引継ぎ中の動作モードの指示を行なう。

【0057】上記に示したマネージメントツールの役割は記憶制御装置 40 及び記憶制御装置 40 に搭載されるパネル 36 を用いて実現することも可能である。この場合はオペレータが記憶制御装置 40 へ引継ぎをせたいロジカルユニットをパネル 36 を用いて入力し、当該ロジカルユニットを担当している記憶制御装置 30 を Inquiry コマンドを用いて探し出す。そして、引継ぎ対象のロジカルユニットを所有している記憶制御装置 30 の AL-PA を記憶しておけば良い。また、引継ぎを行なうロジカルユニット及び当該ロジカルユニットを担当していた記憶制御装置 30 のロジカルユニット引継ぎ中の動作モードはオペレータによってパネル 36 を用いて入力される。

【0058】次に、引継ぎを行なうロジカルユニットが決定した記憶制御装置 40 はロジカルユニット引継ぎ処理に入る。まず、記憶制御装置 40 は記憶制御装置 30 に対して FCP_CMND の 1 種である Mode Sense コマンドを発行し、当該ロジカルユニットの状態を調べる。当該ロジカルユニットの状態が正常であれば引継ぎ処理に入るが、ロジカルユニットが使用不可能な閉塞状態であった場合はマネージメントツールを持つ上位装置 B (50) またはパネル 36 を使ってオペレータへ通知する。

【0059】次に、当該ロジカルユニットを担当していた記憶制御装置 30 のロジカルユニット引継ぎ中の動作モードの設定方法を説明する。これは記憶制御装置 40 が FCP_CMND の 1 種である Mode Select コマンドを発行することで、記憶制御装置 30 の設定変更を行なう。これ以降、ロジカルユニット引継ぎ中に上位装置が記憶制御装置 30 に対して行なう当該ロジカルユニットへの処理要求に対して、記憶制御装置 30 は FCP_REP100 にて Busy ステータスを返す。ここで、上位装置からのコマンド処理要求に対して Busy ステータスで応答し続けた場合、上位装置のコマンド処理要求がタイムアウトになる恐れがあるため、ロジカルユニット引継ぎ処理の時間はタイムアウトとならない範囲で行なわせる。

【0060】次に、記憶制御装置 30 は当該ロジカルユニットに対して、上位装置から処理要求を受領したがまだ処理を行なっていないコマンド (キューイングコマンド) があった場合はそのコマンドの処理を実行し、さらに、当該キャッシュ上に残っていた当該ロジカルユニットに対する全ての Write データをディスクドライブ装置 20 へ書き込む。これにより、ロジカルユニットの担当が記憶制御装置 30 から記憶制御装置 40 へ移った場合でもディスクドライブ装置 20 内のデータの整合性が保証される。

【0061】次に、記憶制御装置 40 は記憶制御装置 3

0から記憶制御装置40へ当該ロジカルユニット引継ぎ時に必要な物理ディスクドライブ制御テーブル110及びロジカルユニット制御テーブル120の格納された制御メモリ32の内容を記憶制御装置30からコピーする。この制御メモリ32のコピーは記憶制御装置40が記憶制御装置30へFCP_CMND90の1種であるMode Senseコマンドを発行し、記憶制御装置30の制御メモリ32の情報を読み出し、その内容を記憶制御装置40の制御メモリ32へ書き込むことで実現する。

【0062】記憶制御装置40はロジカルユニットの引継ぎ時に必要な物理ディスクドライブ制御テーブル110及びロジカルユニット制御テーブル120の情報を受領すると、コマンド処理可能な状態にするために内部処理を実行し、上位装置10からのコマンド処理要求に備える。上位装置10からは当該ロジカルユニットのアクセス経路を既に存在していた記憶制御装置30から新たに追加された記憶制御装置40へ変更することで継続して処理を行なうことができる。さらに、上記で示した記憶制御装置30、40間でロジカルユニットの制御情報をオンライン中に引継ぐ手法は、コンピュータシステムのオンライン中にファイバチャネルネットワーク上に存在する記憶制御装置が同一ファイバチャネルネットワークに接続された別の記憶制御装置からロジカルユニット単位で制御情報を引継ぎ、これ以降の上位装置から当該ロジカルユニットに対して行われるコマンド処理要求を担当し、ファイバチャネルネットワーク上に存在する記憶制御装置間の負荷分散及び複数からなる記憶制御装置が行なっていたいくつかの処理を任意の記憶制御装置へ統合する手段にも適用可能である。

【0063】次に、ロジカルユニットの担当が記憶制御装置30から記憶制御装置40へ移った場合に、オペレータによる上位装置からのアクセス経路の変更を不要とする方法を説明する。図17に示すように、記憶制御装置は予め担当しているロジカルユニットと同数のファイバチャネルアダプタを搭載し、各々にAL-PAを設定しておく。各々のAL-PAは担当するロジカルユニットと組みにして管理する。ロジカルユニットの担当が記憶制御装置30から記憶制御装置40へ移動した場合、ロジカルユニットの制御情報だけでなく、ロジカルユニットと組になっているAL-PAも引継がせる。これにより、上位装置からはロジカルユニットの担当が記憶制御装置30から記憶制御装置40へ移ったこと認識しないため、当該ロジカルユニット引継ぎ以前と同様に処理要求を行なうことができ、アクセス経路変更は不要となる。

【0064】このように、本発明を応用することにより、ファイバチャネルネットワーク60へ記憶制御装置40をオンライン中に追加し、記憶制御装置間の負荷分散を行うことが可能である。さらに同手法はファイバ

チャネルネットワーク60上の複数の記憶制御装置が担当するロジカルユニットを自由に替えていくことが可能であり、既存の記憶制御装置間での負荷分散及び統合が実現できる。

【0065】

【発明の効果】本発明は、ファイバチャネルネットワークに接続された既存の記憶制御装置に対し、当該ファイバチャネルネットワーク上に新たな記憶制御装置をオンライン中に追加し、ロジカルユニットの処理担当を引継がせることで、記憶制御装置間の負荷分散、処理の統合を行うことができ、データ格納システム全体の処理能力の向上が図れる。

【0066】具体的には、ANSI X3T11で標準化されたファイバチャネルネットワークに接続された既存の記憶制御装置が処理能力不足に陥った場合に、新たな記憶制御装置をコンピュータシステムのオンライン中に追加でき、負荷分散が図れる。

【0067】また、本発明は、FC-AL (Fibre Channel Arbitrated Loop) トポロジで最大127台、Fabricトポロジで約1700万台の装置を接続でき、さらにオンライン中にファイバチャネルネットワークへ装置を物理的に接続し、装置側からリンクリセットを発行することにより自動的にファイバチャネルネットワーク内へのログインが可能であるという特長をもつファイバチャネルをインターフェースとして採用することで、2台の記憶制御装置のみならず、オンライン中に複数台の記憶制御装置の追加が可能であり、初期設定にとらわれない柔軟なデータ格納システムの拡張を行なうことができる。

30 【図面の簡単な説明】

【図1】本発明の一実施形態を示すハードウェア構成図である。

【図2】実施形態における記憶制御装置の内部ハードウェア構成図である。

【図3】ファイバチャネルネットワークにFabricトポロジを採用したハードウェア構成図である。

【図4】ファイバチャネルネットワークに複数のFC-ALループを用いたハードウェア構成図である。

【図5】実施形態におけるファイバチャネルがデータをやりとりする基本単位であるフレームのフォーマットを示した図である。

【図6】図5で示したフレームを構成するフレームヘッダのフォーマットを示した図である。

【図7】図5で示したフレームの一構成要素であるFCP_CMNDのペイロードのフォーマット及び当該ペイロードを構成するFCP_CDBのフォーマットを示した図である。

【図8】図5で示したフレームの一構成要素であるFCP_RSPのペイロードのフォーマットを示した図である。

【図9】 Inquiryコマンドのシーケンスを示した図である。

【図10】 Mode Senseコマンドのシーケンスを示した図である。

【図11】 Mode Selectコマンドのシーケンスを示した図である。

【図12】 実施形態における記憶制御装置内部の制御メモリに格納される物理ディスクドライブ制御テーブルのフォーマットを示した図である。

【図13】 実施形態における記憶制御装置内部の制御メモリに格納されるロジカルユニット制御テーブルのフォーマットを示した図である。

【図14】 実施形態におけるファイバチャネルネットワークに新規に追加した記憶制御装置の動作シーケンスを示すフローチャートである

【図15】 実施形態における上位装置に搭載されるマネージメントツールの動作シーケンスを示すフローチャートである

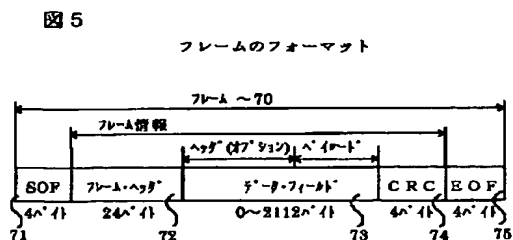
【図16】 実施形態における上位装置に搭載されるマネージメントツールの管理情報を格納するマネージメントツール管理テーブルのフォーマットを示した図である。

【図17】 実施形態において、記憶制御装置間でロジカルユニット引継ぎが行われた後、上位装置からのアクセス経路の変更を不要とすることを実現したデータ格納システム構成図である。

【符号の説明】

10 上位装置A	20 ディスクドライブ装置
21 ロジカルユニット0	22 ロジカルユニット1
30 記憶制御装置	31 ファイバチャネル制御部
32 制御メモリ	33 CPU
34 キャッシュ制御部	35 キャッシュ
36 パネル	40 記憶制御装置 (新規追加)
50 上位装置B (マネージメントツール)	
60 ファイバチャネルネットワーク (FC-AL)	

【図5】



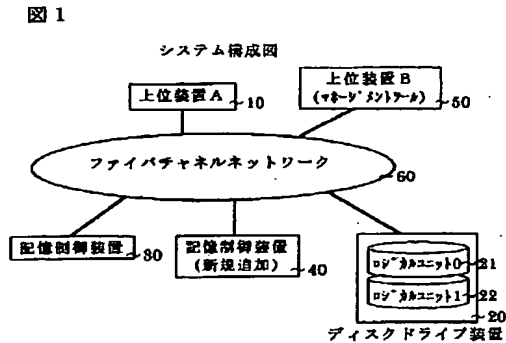
70	フレームのフォーマット
71	SOF (Start Of Frame)
72	フレームヘッダ
73	データフィールド
74	CRC (Cyclic Redundancy Check)
75	EOF (End Of Frame)
80	フレームヘッダのフォーマット
81	D_ID (Destination ID)
82	S_ID (Source ID)
90	FCP_CMNDペイロード (Fibre Channel Protocol for SCSI Command)
91	FCP_LUN (FCP Logical Unit Number)
92	FCP_CNTL (FCP Control)
93	FCP_CDB (FCP Command Descriptor Block)
94	FCP_DL (FCP Data Length)
95	Operation Code
100	FCP_RSPペイロード
101	FCP_STATUS (FCP Status)
110	物理ディスクドライブ制御テーブル
111	物理ディスクドライブ番号
112	物理ディスクドライブ位置
113	記憶容量
114	ブロック数
115	RAIDグループ番号
116	状態
120	ロジカルユニット制御テーブル
121	ロジカルユニット番号
122	RAIDグループ番号
123	RAIDレベル
124	先頭アドレス
125	最終アドレス
130	マネージメントツール管理テーブル
131	記憶制御装置のAL-PA
132	担当LUN

【図8】

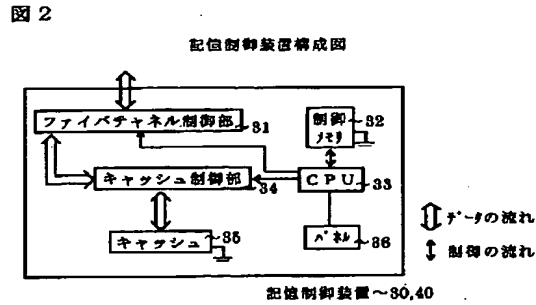
図8 FCP_RSPペイロード100のフォーマット

Bit	31-24	23-16	15-8	7-0
0-1	Reserved			
2	FCP_STATUS (SCSIコマンドの信頼性を設定)			
3	FCP_RESID (未転送データの量を設定)			
4	FCP_SNS_LEN (FCP_SNS_INFO領域で有効なバイト数を設定)			
5	FCP_RSP_LEN (FCP_RSP_INFO領域で有効なバイト数を設定)			
6~n	FCP_RSP_INFO (FCPコマンドの応答情報を設定)			
n-31	FCP_SNS_INFO (SCSIコマンド情報)			

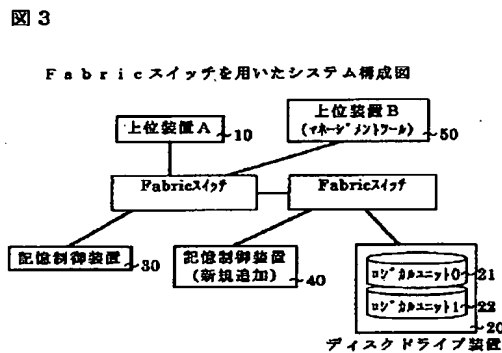
【図1】



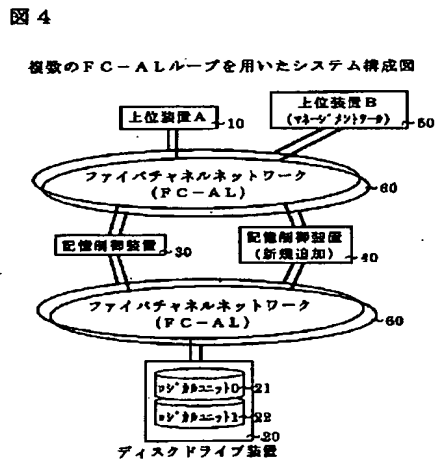
【図2】



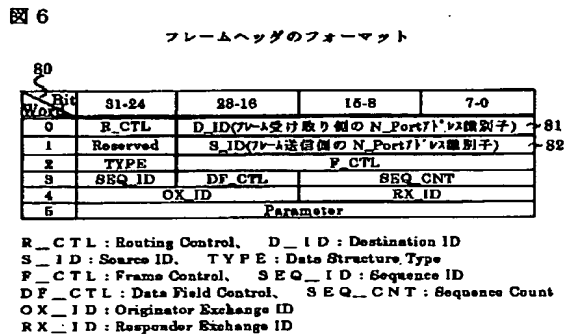
【図3】



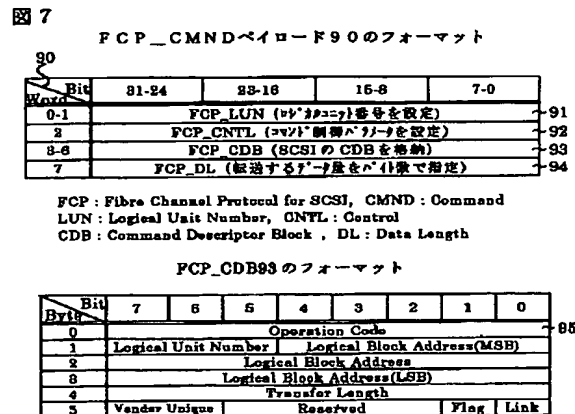
【図4】



【図6】



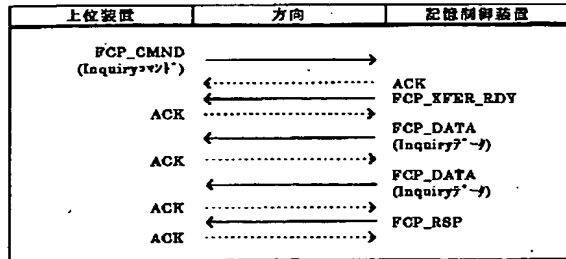
【図7】



【図9】

図9

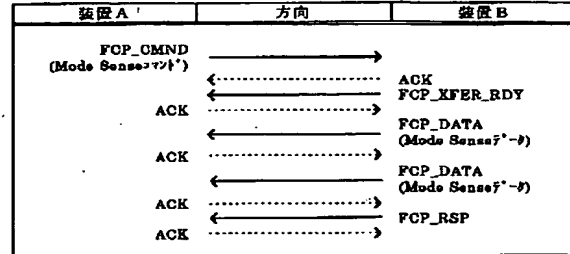
Inquiryコマンドのシーケンス



【図10】

図10

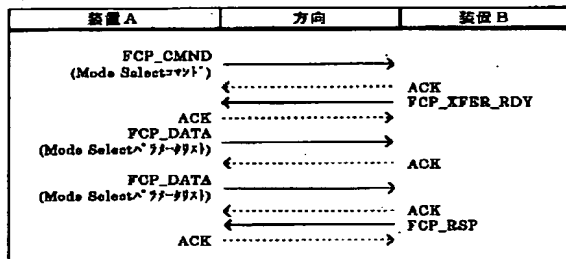
Mode Senseコマンドのシーケンス



【図11】

図11

Mode Selectコマンドのシーケンス



【図12】

図12

物理ディスクドライブ制御テーブル 110

物理ディスク ドライブ番号	物理ディスク ドライブ位置	記憶 容量	ブロック 数	RAIDグループ 番号	状態	記憶装置 の種類
0						
1						
...						

【図16】

図16

マネージメントツール管理テーブル 130

記憶制御装置のAL-PA	担当 LUN

【図13】

図13

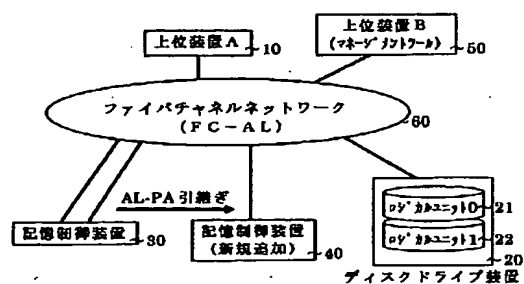
ロジカルユニット制御テーブル 120

ロジカルユニット番号	RAIDグループ 番号	RAIDレベル	先頭アドレス	最終アドレス
0				
1				
...				

【図17】

図17

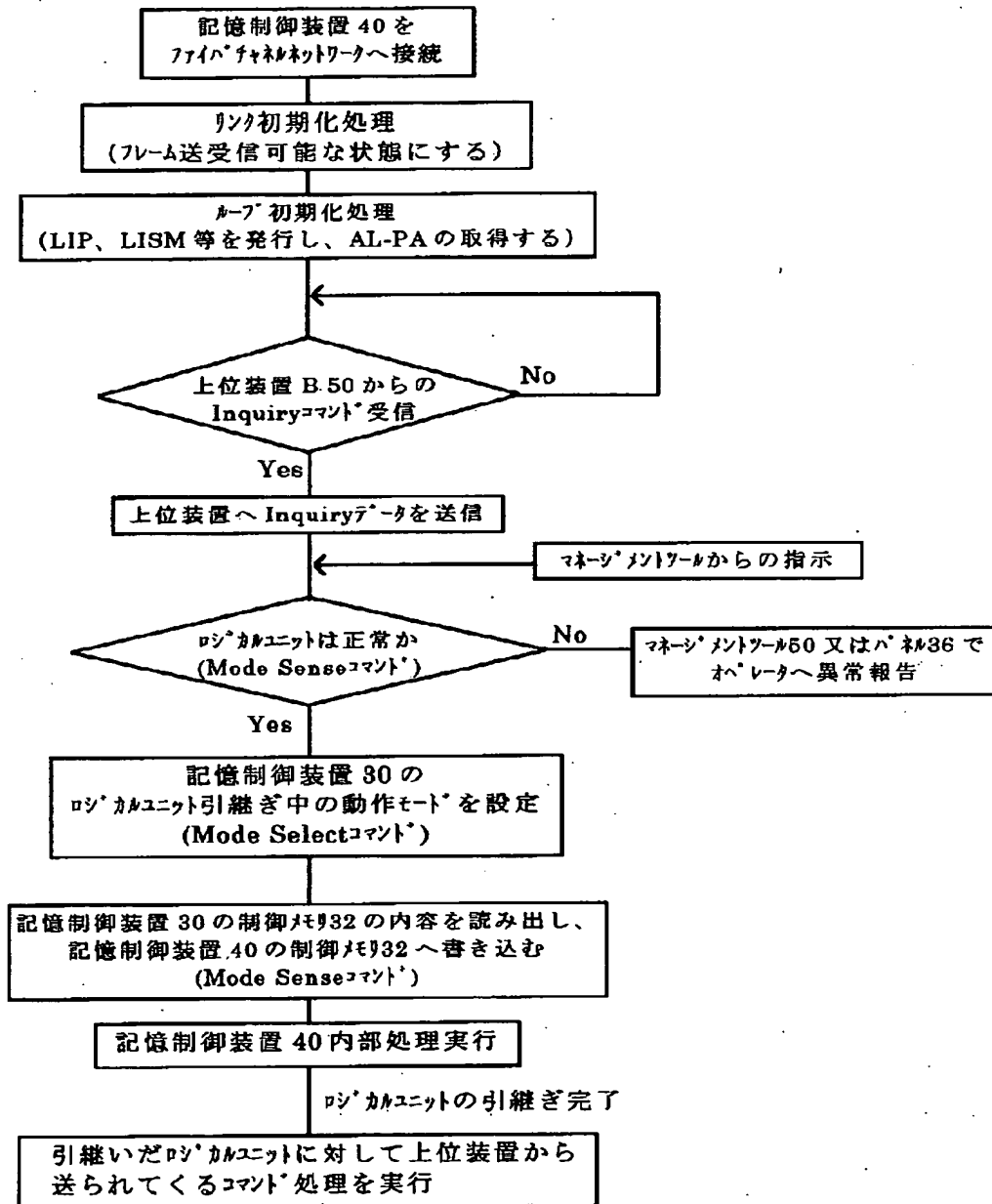
システム構成図



【図14】

図14

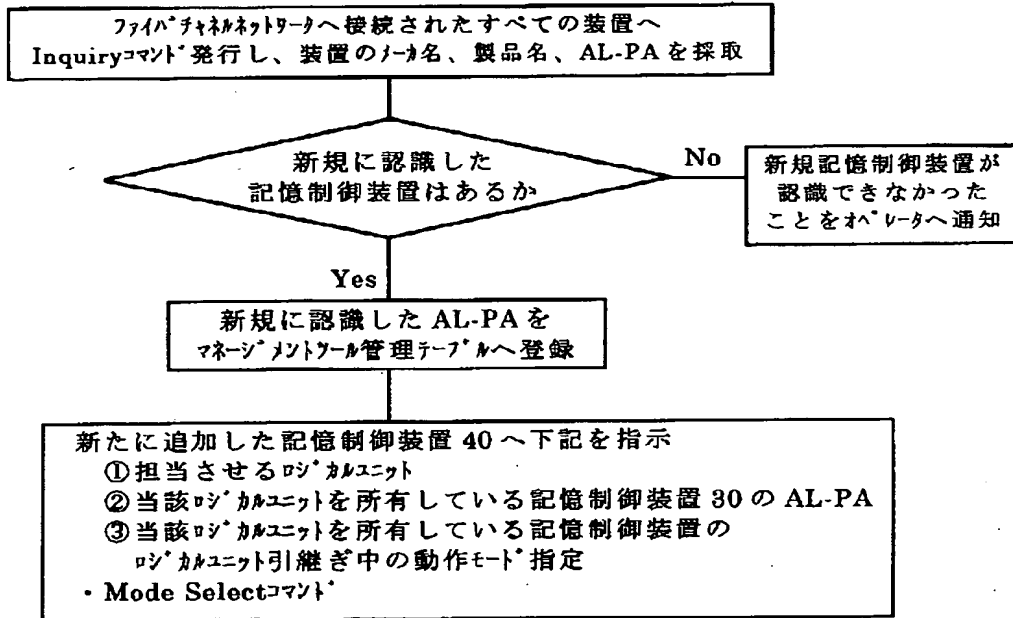
記憶制御装置40の動作シーケンス



【図15】

図15

上位装置B 50に搭載された
マネージメントツールの動作シーケンス



フロントページの続き

(72)発明者 岸本 哲哉
神奈川県小田原市国府津2880番地 株式会
社日立製作所ストレージシステム事業部内
(72)発明者 岩崎 秀彦
神奈川県小田原市国府津2880番地 株式会
社日立製作所ストレージシステム事業部内

(72)発明者 村岡 健司
神奈川県小田原市国府津2880番地 株式会
社日立製作所ストレージシステム事業部内
(72)発明者 高本 賢一
神奈川県小田原市国府津2880番地 株式会
社日立製作所ストレージシステム事業部内

Fターム(参考) 5B065 BA01 BA03 BA04 BA06 BA07
CA11 CA13 CA30 CH13 EA33
ZA03 ZA05 ZA13 ZA14